



UNITED STATES PATENT AND TRADEMARK OFFICE

UNITED STATES DEPARTMENT OF COMMERCE
United States Patent and Trademark Office
Address: COMMISSIONER FOR PATENTS
P.O. Box 1450
Alexandria, Virginia 22313-1450
www.uspto.gov

APPLICATION NO.	FILING DATE	FIRST NAMED INVENTOR	ATTORNEY DOCKET NO.	CONFIRMATION NO.
10/600,475	06/20/2003	Chris J.C. Burges	MCS-018-03	6335
7590 LYON & HARR, L.L.P Suite 800 300 Esplanade Drive Oxnard, CA 93036-1274	04/29/2009		EXAMINER BAKER, MATTHEW H	
			ART UNIT 2626	PAPER NUMBER
			MAIL DATE 04/29/2009	DELIVERY MODE PAPER

Please find below and/or attached an Office communication concerning this application or proceeding.

The time period for reply, if any, is set in the attached communication.

Office Action Summary	Application No.	Applicant(s)	
	10/600,475	BURGES ET AL.	
	Examiner	Art Unit	
	Matthew Baker	2626	

-- The MAILING DATE of this communication appears on the cover sheet with the correspondence address --

Period for Reply

A SHORTENED STATUTORY PERIOD FOR REPLY IS SET TO EXPIRE 3 MONTH(S) OR THIRTY (30) DAYS, WHICHEVER IS LONGER, FROM THE MAILING DATE OF THIS COMMUNICATION.

- Extensions of time may be available under the provisions of 37 CFR 1.136(a). In no event, however, may a reply be timely filed after SIX (6) MONTHS from the mailing date of this communication.
- If NO period for reply is specified above, the maximum statutory period will apply and will expire SIX (6) MONTHS from the mailing date of this communication.
- Failure to reply within the set or extended period for reply will, by statute, cause the application to become ABANDONED (35 U.S.C. § 133). Any reply received by the Office later than three months after the mailing date of this communication, even if timely filed, may reduce any earned patent term adjustment. See 37 CFR 1.704(b).

Status

- 1) Responsive to communication(s) filed on 02-17-2009.
 2a) This action is FINAL. 2b) This action is non-final.
 3) Since this application is in condition for allowance except for formal matters, prosecution as to the merits is closed in accordance with the practice under *Ex parte Quayle*, 1935 C.D. 11, 453 O.G. 213.

Disposition of Claims

- 4) Claim(s) 1,5-16,19,20 and 22-29 is/are pending in the application.
 4a) Of the above claim(s) _____ is/are withdrawn from consideration.
 5) Claim(s) _____ is/are allowed.
 6) Claim(s) 1,5-16,19,20 and 22-29 is/are rejected.
 7) Claim(s) _____ is/are objected to.
 8) Claim(s) _____ are subject to restriction and/or election requirement.

Application Papers

- 9) The specification is objected to by the Examiner.
 10) The drawing(s) filed on 20 June 2003 is/are: a) accepted or b) objected to by the Examiner.
 Applicant may not request that any objection to the drawing(s) be held in abeyance. See 37 CFR 1.85(a).
 Replacement drawing sheet(s) including the correction is required if the drawing(s) is objected to. See 37 CFR 1.121(d).
 11) The oath or declaration is objected to by the Examiner. Note the attached Office Action or form PTO-152.

Priority under 35 U.S.C. § 119

- 12) Acknowledgment is made of a claim for foreign priority under 35 U.S.C. § 119(a)-(d) or (f).
 a) All b) Some * c) None of:
 1. Certified copies of the priority documents have been received.
 2. Certified copies of the priority documents have been received in Application No. _____.
 3. Copies of the certified copies of the priority documents have been received in this National Stage application from the International Bureau (PCT Rule 17.2(a)).

* See the attached detailed Office action for a list of the certified copies not received.

Attachment(s)

- | | |
|--|---|
| 1) <input type="checkbox"/> Notice of References Cited (PTO-892) | 4) <input type="checkbox"/> Interview Summary (PTO-413) |
| 2) <input type="checkbox"/> Notice of Draftsperson's Patent Drawing Review (PTO-948) | Paper No(s)/Mail Date. _____ . |
| 3) <input type="checkbox"/> Information Disclosure Statement(s) (PTO/SB/08) | 5) <input type="checkbox"/> Notice of Informal Patent Application |
| Paper No(s)/Mail Date _____. | 6) <input type="checkbox"/> Other: _____ . |

DETAILED ACTION

Continued Examination Under 37 CFR 1.114

1. A request for continued examination under 37 CFR 1.114, including the fee set forth in 37 CFR 1.17(e), was filed in this application after final rejection. Since this application is eligible for continued examination under 37 CFR 1.114, and the fee set forth in 37 CFR 1.17(e) has been timely paid, the finality of the previous Office action has been withdrawn pursuant to 37 CFR 1.114. Applicant's submission filed on 9 January 2009 has been entered.

Specification

2. The abstract of the disclosure is objected to because the sentence "The anchor model outputs are mapped to frame tags to that all speech corresponding *to* a single frame tag comes from a single speaker." Should read --The anchor model outputs are mapped to frame tags to that all speech corresponding *so* a single frame tag comes from a single speaker. -- Correction is required. See MPEP § 608.01(b).

3. Claims 24-27 and 29 are objected to because of the following informalities: Claim 24 recites "A computer-readable medium having computer-executable instructions for processing audio data..." which should read -- A computer-readable medium encoded with computer-executable instructions for processing audio data, which when executed perform...--. Appropriate correction is required.

Response to Amendment

4. Applicant have filed an amendment 9 February 2009. Claims 1, 5-16, 19, 20, and 22-27 are pending. Applicant has argued to traverse the rejection of the pending claims under 35 USC 103(a).

Response to Arguments

5. Applicant's arguments filed 9 January 2009 have been fully considered but they are not persuasive. The arguments will be addressed in detail below.

Applicant argues that “neither Sturim et al. nor Waibel et al. disclose “obtaining a preliminary output of the plurality of anchor models from the time-delay neural network during training of the TDNN classifiers before final nonlinearities are applied by the second layer in order to generate an output of the plurality of anchor models (Remarks, p. 12).”

Waibel teaches “obtaining a preliminary output of the plurality of anchor models from the time-delay neural network during training of the TDNN classifiers (*Each unit in the first hidden layer now receives input (via 48 weighted connections) from the coefficients in the 3 frame window*, p. 330, ¶ 1) before final nonlinearities are applied by the second layer in order to generate an output of the plurality of anchor models. (*Finally, the output is obtained by integrating (summing) the evidence from each of the 3 units in hidden layer 2 over time and connecting it to its pertinent output unit (shown in Fig. 2 over 9 frames for the “B” output unit).* In practice, this summation is implemented simply as another nonlinear (sigmoid function is applied here as well) TDNN unit which has fixed equal weights to a row of unit firings over time in hidden layer 2, p. 330, ¶ 3).” Applicant should refer to Fig. 2 to further understand that the

“integration” process is the final nonlinearity to be applied (between Hidden Layer 2 and Output Layer). The argument is not convincing.

6. Applicant argues that Hermansky does not teach the newly added limitation of “obtaining a preliminary output of the plurality of anchor models from the time-delay neural network during training of the TDNN classifiers before final nonlinearities are applied by the second layer...” (Remarks, p. 12). This argument is moot as Waibel has been cited to teach said limitation above.

7. Applicant argues that the combination of Sutirm in view of Waibel in few of Hermansky in further view of Lavagetto is improper because the limitations of the independent claims are not taught (Remarks, p. 14). This argument is moot as Waibel has been cited to teach said limitation above.

8. Applicant argues that the combination of Sutirm in view of Waibel in few of Hermansky in further view of Liu is improper because the limitations of the independent claims are not taught. This argument is moot as Waibel has been cited to teach said limitation above.

Claim Rejections - 35 USC § 103

9. The text of those sections of Title 35, U.S. Code not included in this action can be found in a prior Office action.

10. Claims 1, 6-10, 12-14, 16, 20, and 22-27 are rejected under 35 U.S.C. 103(a) as being unpatentable over Sturim (“Speaker Indexing in Large Audio Databases Using Anchor Models” 2001) in view of Waibel (“Phoneme Recognition Using Time-Delay Neural Networks” IEEE 1989), further in view of Hermansky (US Patent 7,254,538).

11. As per claim 1, ***Sturim*** discloses a method for processing audio data, comprising:

applying the plurality of anchor models to the audio data (Section 1. Introduction, *a target utterance is characterized using anchor models derived from a predetermined set of speakers*);

normalizing the modified feature vector output to generate normalized anchor model output (section 2. Anchor Models, second paragraph, *each anchor model yields a likelihood score, where the combination of scores are used to form a N-dimensional characterization vector*, and the fifth paragraph to the sixth paragraph, *a comparison is done between normalized data and non-normalized output data, therefore normalization must have been done*);

mapping the normalized output of the plurality of anchor models into frame tags and producing the frame tags (section 2. Anchor Models, *speaker characterization vectors are mapped onto a speaker space. A determination of the speaker is made based on the location of the vector within speaker space*).

Sturim does not disclose training time-delay neural network (TDNN) classifiers using a time-delay neural network that uses first layer followed by a second layer having a nonlinearity, and using discriminatively trained classifiers that are time-delay neural network (TDNN) classifiers to produce a plurality of anchor model output models. However, **Sturim** does disclose that anchor models, previously trained during a training phrase, can consist of any method of speech representation (section 1. Introduction and section 2. Anchor Models, first paragraph). In addition, **Waibel** discloses a speech processing system that uses a TDNN for speech representation (Abstract), where a TDNN is a type of convolutional neural network. The TDNN of **Waibel** discloses obtaining a preliminary output of the plurality of anchor models from the

time-delay neural network during training of the TDNN classifiers (*Each unit in the first hidden layer now receives input (via 48 weighted connections) from the coefficients in the 3 frame window*, p. 330, ¶ 1) before final nonlinearities are applied by the second layer (*Finally, the output is obtained by integrating (summing) the evidence from each of the 3 units in hidden layer 2 over time and connecting it to its pertinent output unit (shown in Fig. 2 over 9 frames for the “B” output unit). In practice, this summation is implemented simply as another nonlinear (sigmoid function is applied here as well) TDNN unit which has fixed equal weights to a row of unit firings over time in hidden layer 2*, p. 330, ¶ 3).

Therefore it would have been obvious to one of ordinary skill in the art at the time of the invention to use a discriminatively-trained classifiers of a convolutional neural network that was previously trained using a training technique that included non-linear processing step in **Sturim**, since the time-delay structure enables the system to discover the temporal relationship among acoustic features independent of the position in time, as indicated in **Waibel** (Abstract).

Sturim also does not disclose obtaining a preliminary output of the plurality of anchor models from the time-delay neural network during training of the TDNN classifiers before final nonlinearities are applied by the second layer in order to generate an output of the plurality of anchor models. **Hermansky** discloses a system where the final nonlinearity is omitted (column 2 lines 23-25 and column 3 lines 32-35). **Hermansky** discloses a system where a neural network is combined with HMM (Gaussian mixture model) models for speech recognition, and the output of the neural network is adjusted prior to being input into the HMM. The final nonlinearity of the neural network is omitted, thus adjusting the posterior probabilities to make it more clearly Gaussian and optimizing it for a HMM system.

Therefore it would have been obvious to one of ordinary skill in the art at the time of the invention to obtaining a preliminary output of the plurality of anchor models from a time-delay neural network before the second layer having a nonlinearity is applied in order to generate an output of the plurality of anchor models in **Sturim** and **Waibel**, since one of ordinary skill has good reason to pursue the options within his or her technical grasp in order to achieve the predictable result of manipulating the output of a neural network, thus optimizing it for further processing.

As per claim 20, **Sturim** discloses a method for processing audio data containing a plurality of speakers, comprising:

applying the plurality of anchor models to the audio data (Section 1. Introduction, *a target utterance is characterized using anchor models derived from a predetermined set of speakers*);

normalizing the modified feature vector output to generate normalized anchor model output (section 2. Anchor Models, second paragraph, *each anchor model yields a likelihood score, where the combination of scores are used to form a N-dimensional characterization vector*, and the fifth paragraph to the sixth paragraph, *a comparison is done between normalized data and non-normalized output data, therefore normalization must have been done*);

mapping the normalized output of the plurality of anchor models into frame tags and producing the frame tags (section 2. Anchor Models, *speaker characterization vectors are mapped onto a speaker space. A determination of the speaker is made based on the location of the vector within speaker space*).

Wherein discriminatively-trained classifiers were previously trained using a training set containing a set of training speakers, and wherein the plurality of speakers is not in the set of training speakers ((section 1. Introduction, *speakers of the target utterance are not members of the training set*).

Sturim does not disclose training time-delay neural network (TDNN) classifiers using a time-delay neural network that uses a first layer followed by a second layer having a nonlinearity, and using discriminatively trained classifiers that are time-delay neural network (TDNN) classifiers to produce a plurality of anchor model output models. However, **Sturim** does disclose that anchor models, previously trained during a training phrase, can consist of any method of speech representation (section 1. Introduction and section 2. Anchor Models, first paragraph). In addition, **Waibel** discloses a speech processing system that uses a TDNN for speech representation (Abstract), where a TDNN is a type of convolutional neural network. The TDNN of **Waibel** discloses obtaining a preliminary output of the plurality of anchor models from the time-delay neural network during training of the TDNN classifiers (*Each unit in the first hidden layer now receives input (via 48 weighted connections) from the coefficients in the 3 frame window*, p. 330, ¶ 1) before final nonlinearities are applied by the second layer (*Finally, the output is obtained by integrating (summing) the evidence from each of the 3 units in hidden layer 2 over time and connecting it to its pertinent output unit (shown in Fig. 2 over 9 frames for the “B” output unit). In practice, this summation is implemented simply as another nonlinear (sigmoid function is applied here as well) TDNN unit which has fixed equal weights to a row of unit firings over time in hidden layer 2*, p. 330, ¶ 3).

Therefore it would have been obvious to one of ordinary skill in the art at the time of the invention to use a discriminatively-trained classifiers of a convolutional neural network that was previously trained using a training technique that included non-linear processing step in **Sturim**, since the time-delay structure enables the system to discover the temporal relationship among acoustic features independent of the position in time, as indicated in **Waibel** (Abstract).

Sturim also does not disclose obtaining a preliminary output of the plurality of anchor models from a time-delay neural network before the second layer having a nonlinearity is applied in order to generate an output of the plurality of anchor models. **Hermansky** discloses a system where the final nonlinearity is omitted (column 2 lines 23-25 and column 3 lines 32-35).

Hermansky discloses a system where a neural network is combined with HMM (Gaussian mixture model) models for speech recognition, and the output of the neural network is adjusted prior to being input into the HMM. The final nonlinearity of the neural network is omitted, thus adjusting the posterior probabilities to make it more clearly Gaussian and optimizing it for a HMM system.

Therefore it would have been obvious to one of ordinary skill in the art at the time of the invention to obtaining a preliminary output of the plurality of anchor models from a time-delay neural network before the second layer having a nonlinearity is applied in order to generate an output of the plurality of anchor models in **Sturim** and **Waibel**, since one of ordinary skill has good reason to pursue the options within his or her technical grasp in order to achieve the predictable result of manipulating the output of a neural network, making it optimized for further processing.

Sturim also does not explicitly state constructing a list of start and stop times for each of the plurality of speakers based on the frame tags, nor using discriminatively-trained classifiers that are time-delay neural network TDNN) classifiers to produce a plurality of anchor model outputs. However, **Sturim** does disclose a system that can be used to retrieve messages from an archive (section 4. Speaker Indexing). In order to retrieve the messages, the system must know where the messages start and stop, and therefore must determine start and stop times for each speaker. **Sturim** also discloses that anchor models, previously trained during a training phrase, can consist of any method of speech representation (section 1. Introduction and section 2. Anchor Models, first paragraph). In addition, **Waibel** discloses a speech processing system that uses a TDNN for the speech representation. (Abstract). **Sturim** and **Waibel** both disclose system for improved speech processing, and are therefore analogous art.

Therefore it would have been obvious to one of ordinary skill in the art at the time of the invention to construct a list of start and stop times in **Sturim**, since start and stop times can be used to reliably retrieve and playback saved messages corresponding to a specific speaker.

It would also have been obvious to one of ordinary skill in the art at the time of the invention to have a TDNN as an anchor model in **Sturim**, since the time-delay structure enables the system to discover the temporal relationship among acoustic features independent of the position in time, as indicated in **Waibel** (Abstract).

As per claim 6, **Sturim** in view of **Waibel** disclose the method as set forth in claim 1, and **Waibel** further discloses pre-processing the audio data to generate input feature vectors for the

discriminatively-trained classifier (section II A, *melscale spectral coefficients are derived from the input speech, then input to the network*).

Therefore it would have been obvious to one of ordinary skill in the art at the time of the invention to pre-process the audio data to generate input feature vectors in **Sturim**, since it would provide a reliable set of feature vectors, which can be easily applied to the classifier for further processing.

As per claims 7, 8 and 22, **Sturim** in view of **Waibel** disclose the method as set forth in claims 1 and 20 and **Sturim** further discloses normalizing a feature vector output of the discriminatively-trained classifier (section 2. Anchor Models, second paragraph, *each anchor model yields a likelihood score, where the combination of scores are used to form a N-dimensional characterization vector*, and the fifth paragraph to the sixth paragraph, *a comparison is done between normalized data and non-normalized output data, therefore normalization must have been done*). **Sturim** does not explicitly state wherein the normalized feature vectors are vectors of unit length. However Official notice is taken that it is old and well known to normalize a vector to a vector of unit length. During vector processing, feature vectors are often normalized to a unit vector for simplicity of computation.

Therefore it would have been obvious to one of ordinary skill in the art at the time of the invention to normalize a vector output of the classifier to a unit vector in **Sturim**, since it would produce a simplified feature vector, enabling simplified processing which then reserves computational resources.

As per claim 9, **Sturim** in view of **Waibel** disclose the method as set forth in claim 1, however **Sturim** does not explicitly disclose accepting a plurality of input feature vectors corresponding to audio features contained in the audio data, and applying the discriminatively-trained classifier to the plurality of input feature vectors to produce a plurality of anchor model outputs. However, **Sturim** does disclose applying input data to a trained anchor models to produce anchor model outputs (section 2. Anchor Models). In addition, **Waibel** discloses accepting a plurality of feature vectors corresponding to audio features contained in the audio data (section II A, *melscale spectral coefficients are derived from the input speech, then input to the network*), and applying the discriminatively-trained classifier to the plurality of input feature vectors to produce a plurality of model outputs (Abstract, *a TDNN is used for speech processing*).

Therefore it would have been obvious to one of ordinary skill in the art at the time of the invention to accept a plurality of input feature vectors corresponding to audio features contained in the audio data, and apply the discriminatively-trained classifier to the plurality of input feature vectors to produce a plurality of anchor model outputs in **Sturim**, since it would provide a reliable set of feature vectors which can be easily applied to a TDNN, where the time-delay structure enables the system to discover the temporal relationship among acoustic features independent of the position in time, as indicated in **Waibel** (Abstract).

As per claim 10, **Sturim** in view of **Waibel** disclose the method as set forth in claim 1, and **Sturim** further discloses wherein the mapping comprises clustering anchor model outputs from

the discriminatively-trained classifier into separate clusters using a clustering technique, and associating a frame tag to each separate cluster (section 2. Anchor Models, *speaker characterization vectors are mapped onto a speaker space. A determination of the speaker is made based on the location of the vector within speaker space*).

As per claim 12, **Sturim** in view of **Waibel** disclose the method as set forth in claim 1, and **Sturim** further discloses training the discriminatively-trained classifier using a speaker training set containing a plurality of known speakers (section 1. Introduction, *anchor models are derived from a set of predetermined speakers*). **Sturim** does not explicitly disclose pre-processing the speaker training set and the audio data in the same manner to provide a consistent input to the discriminatively trained classifier. However, **Waibel** discloses pre-processing of audio data (section II A, *melscale spectral coefficients are derived from the input speech, then input to the network*).

Therefore it would have been obvious to one of ordinary skill in the art at the time of the invention to pre-process the speaker training set and the audio data in the same manner in **Sturim**, since it would provide reliable data input to the classifier, which would then provide a reliable and useful result.

As per claims 13 and 23, neither **Sturim** in view of **Waibel** explicitly disclose computer-readable medium having computer-executable instructions for performing the method recited in claims 1 and 20. However, the method of **Sturim** requires considerable computation and processing, and modern computer systems can perform the same computations considerably faster, and with

higher accuracy, than any human would. In addition, **Waibel** states that the disclosed system was implemented using C and Fortran (page 331, second column), both common programming languages used to execute computer readable instructions.

Therefore it would have been obvious to perform the method of claim 1 on a computer-readable medium in **Sturim**, since a computer would enable faster processing, saving time and assuring accuracy.

As per claim 14, **Sturim** discloses a computer-implemented process for processing audio data, comprising:

applying a plurality of anchor models to the audio data (Section 1. Introduction, *a target utterance is characterized using anchor models derived from a predetermined set of speakers*);

normalizing the modified feature vector output to generate normalized anchor model output (section 2. Anchor Models, second paragraph, *each anchor model yields a likelihood score, where the combination of scores are used to form a N-dimensional characterization vector*, and the fifth paragraph to the sixth paragraph, *a comparison is done between normalized data and non-normalized output data, therefore normalization must have been done*);

mapping the normalized anchor model output into frame tags and producing frame tags (section 2. Anchor Models, *speaker characterization vectors are mapped onto a speaker space. A determination of the speaker is made based on the location of the vector within speaker space*).

Sturim does not disclose the plurality of anchor models comprising discriminatively-trained classifiers of a convolutional neural network that were previously trained using a training technique that uses a first layer followed by a second layer having a nonlinearity, and obtaining a preliminary output of the plurality of anchor models from the convolutional neural network before the second layer having the nonlinearity is applied in order to generate a modified feature vector output. However, **Sturim** does disclose that anchor models, previously trained during a training phrase, can consist of any method of speech representation (section 1. Introduction and section 2. Anchor Models, first paragraph). In addition, **Waibel** discloses a speech processing system that uses a TDNN for speech representation (Abstract), where a TDNN is a type of convolutional neural network. The TDNN of **Waibel** discloses obtaining a preliminary output of the plurality of anchor models from the time-delay neural network during training of the TDNN classifiers (*Each unit in the first hidden layer now receives input (via 48 weighted connections) from the coefficients in the 3 frame window*, p. 330, ¶ 1) before final nonlinearities are applied by the second layer (*Finally, the output is obtained by integrating (summing) the evidence from each of the 3 units in hidden layer 2 over time and connecting it to its pertinent output unit (shown in Fig. 2 over 9 frames for the “B” output unit). In practice, this summation is implemented simply as another nonlinear (sigmoid function is applied here as well) TDNN unit which has fixed equal weights to a row of unit firings over time in hidden layer 2*, p. 330, ¶ 3).

Hermansky discloses a system where the final nonlinearity is omitted (column 2 lines 23-25 and column 3 lines 32-35). **Hermansky** discloses a system where a neural network is combined with HMM (Gaussian mixture model) models for speech recognition, and the output of the neural network is adjusted prior to being input into the HMM. The final nonlinearity of the neural

network is omitted, thus adjusting the posterior probabilities to make it more clearly Gaussian and optimizing it for a HMM system.

Therefore it would have been obvious to one of ordinary skill in the art at the time of the invention to use a discriminatively-trained classifiers of a convolutional neural network that was previously trained using a training technique that included non-linear processing step in *Sturim*, since the time-delay structure enables the system to discover the temporal relationship among acoustic features independent of the position in time, as indicated in *Waibel* (Abstract).

Therefore it would also have been obvious to one of ordinary skill in the art at the time of the invention to obtaining a preliminary output of the plurality of anchor models from a time-delay neural network before the second layer having a nonlinearity is applied in order to generate an output of the plurality of anchor models in *Sturim* and *Waibel*, since one of ordinary skill has good reason to pursue the options within his or her technical grasp in order to achieve the predictable result of manipulating the output of a neural network, making it optimized for further processing.

As per claim 16, *Sturim* in view of *Waibel* disclose the system of claim 14, and *Waibel* further discloses wherein the training technique employs a mean-square error metric (section II B, first paragraph). *Waibel* also discloses that there are many learning techniques for the optimization of neural networks, including mean-squared error.

Therefore it would have been obvious to one of ordinary skill in the art at the time of the invention to use the mean-square error during training in *Sturim*, since one of ordinary skill in

the art at the time of the invention has good reason to pursue the options within his to her technical grasp.

As per claim 19, this claim recites limitations similar to claim 8, and is therefore rejected for similar reasons.

As per claim 24, **Sturim** disclose a computer-readable medium having computer-executable instructions for processing audio data, comprising:

training anchor models to be used to produce anchor models outputs, and (Section 1. Introduction, *a target utterance is characterized using anchor models derived from a predetermined set of speakers* and section 2. Anchor Models, *speaker characterization vectors are mapped onto a speaker space. A determination of the speaker is made based on the location of the vector within speaker space*).

normalizing the modified plurality of anchor model output to generate normalized anchor model outputs (section 2. Anchor Models, second paragraph, *each anchor model yields a likelihood score, where the combination of scores are used to form a N-dimensional characterization vector*, and the fifth paragraph to the sixth paragraph, *a comparison is done between normalized data and non-normalized output data, therefore normalization must have been done*);

clustering anchor model outputs into frame tags of speakers (Section 1. Introduction, *a target utterance is characterized using anchor models derived from a predetermined set of speakers* and section 2. Anchor Models, *speaker characterization vectors are mapped onto a speaker space. A determination of the speaker is made based on the location of the vector within speaker space*).

Sturim does not disclose training a discriminatively-trained classifiers that are time-delay neural network (TDNN) classifiers in a discriminative manner on a convolutional neural network using a training technique such that the training occurs during a training phase to generate parameters that can be used at a later time by the TDNN classifiers and includes two layers with a first layer including a one-dimensional convolution followed by a second layer having a nonlinearity, and using the TDNN classifiers to produce a plurality of anchor model outputs. However, **Sturim** does disclose that anchor models, previously trained during a training phrase, can consist of any method of speech representation (section 1. Introduction and section 2. Anchor Models, first paragraph). In addition, **Waibel** discloses a speech processing system that uses a TDNN for speech representation (Abstract), where a TDNN is a type of convolutional neural network. The TDNN of **Waibel** discloses obtaining a preliminary output of the plurality of anchor models from the time-delay neural network during training of the TDNN classifiers (*Each unit in the first hidden layer now receives input (via 48 weighted connections) from the coefficients in the 3 frame window, p. 330, ¶ 1*) before final nonlinearities are applied by the second layer (*Finally, the output is obtained by integrating (summing) the evidence from each of the 3 units in hidden layer 2 over time and connecting it to its pertinent output unit (shown in Fig. 2 over 9 frames for the “B” output unit)*). In practice, this summation is implemented simply

as another nonlinear (sigmoid function is applied here as well) TDNN unit which has fixed equal weights to a row of unit firings over time in hidden layer 2, p. 330, ¶ 3).

Therefore it would have been obvious to one of ordinary skill in the art at the time of the invention to train a discriminatively-trained classifier that is a time-delay neural network (TDNN) in a discriminative manner on a convolutional neural network using a training technique that includes a non-linear processing step such that the training occurs during a training phase to generate parameters that can be used at a later time by the TDNN classifier, using the discriminatively-trained classifiers that are time-delay neural network (TDNN) classifiers to produce a plurality of anchor model outputs in **Sturim**, since the time-delay structure enables the system to discover the temporal relationship among acoustic features independent of the position in time, as indicated in **Waibel** (Abstract).

Neither **Sturim** nor **Waibel** disclose obtaining during training the plurality of anchor model outputs from the convolutional neural network prior to application of final nonlinearities by the second layer to generate a modified plurality of anchor model outputs. However, **Sturim** disclose the use of anchor models, where anchor models are trained classifiers and the output is input into another machine-learning algorithm, such as a clustering algorithm. In addition, **Hermansky** discloses a system where the final nonlinearity is omitted (column 2 lines 23-25 and column 3 lines 32-35). **Hermansky** discloses a system where a neural network is combined with HMM (Gaussian mixture model) models for speech recognition, and the output of the neural network is adjusted prior to being input into the HMM. The final nonlinearity of the neural network is omitted, thus adjusting the posterior probabilities to make it more clearly Gaussian and optimizing it for a HMM system.

Therefore it would have been obvious to one of ordinary skill in the art at the time of the invention to obtain during training a preliminary output of the plurality of anchor models from a time-delay neural network before the second layer having a nonlinearity is applied in order to generate an output of the plurality of anchor models in **Sturim** and **Waibel**, since one of ordinary skill has good reason to pursue the options within his or her technical grasp in order to achieve the predictable result of manipulating the output of a neural network, making it optimized for further processing.

As per claim 25, **Sturim** in view of **Waibel** disclose the computer-readable medium of claim 24, and **Waibel** further discloses pre-processing a speaker training set during the training phase to produce a first set of input feature vectors for the discriminatively-trained classifier (section II A, *melscale spectral coefficients are derived from the input speech, then input to the network*). However, neither **Sturim** nor **Waibel** disclose pre-processing a speaker training set during a validation phase to produce a first set of input feature vectors for the discriminatively-trained classifier. However, by applicants own admission (specification page 17, second paragraph) validation sets are old and well known.

Therefore it would have been obvious to one of ordinary skill in the art at the time of the invention to pre-process the audio data to generate input feature vectors in a training and validation phase in **Sturim**, since it would provide a reliable set of feature vectors, which can be easily applied to the classifier for further processing.

As per claim 26, **Sturim** in view of **Waibel** disclose the computer-readable medium of claim 25, however **Sturim** does not explicitly disclose pre-processing the audio data during the use phase to produce a second set of input feature vectors for the discriminatively-trained classifier, the pre-processing of the audio data being performed in the same manner as the pre-processing of the speaker training set. However, **Waibel** discloses pre-processing of audio data (However, **Waibel** discloses pre-processing of audio data (*section II A, melscale spectral coefficients are derived from the input speech, then input to the network*)).

Therefore it would have been obvious to one of ordinary skill in the art at the time of the invention to pre-process the speaker training set and the audio data in the same manner in **Sturim**, since it would provide reliable data input to the classifier, which would provide a reliable and useful result.

As per claim 27, this claim recites limitations similar to those recited in claim 8, and is therefore rejected for similar reasons.

Claims 5,15 are rejected under 35 U.S.C. 103(a) as being unpatentable over **Sturim** in view of **Waibel**, further in view of **Hermansky** as applied to claims 1 and 14 above, and further in view of **Lavagetto** (“Time-Delay Neural Network for Estimating Lip Movements from Speech Analysis: A useful Tool in Audio-Video Synchronization” IEEE 1997).

Sturim in view of **Waibel** disclose the method as set forth in claim 1 and 14, however neither explicitly disclose further training the TDNN classifier using cross entropy. However, by

applicant's own admission training using cross entropy is well known (specification page 29). In addition, *Lavagetto* discloses that training a time-delay neural network can be done with either cross entropy or mean-square error (page 789-790).

Therefore it would have been obvious to one of ordinary skill in the art at the time of the invention to use cross entropy or mean-square error to train the TDNN in *Sturim* and *Waibel*, since cross entropy and mean-square error provide figures for validating estimates provided by each network independent from the network structure itself, as indicated in *Lavagetto* (page 789, section IV. Learning Criteria for TDNN Training).

Claims 11 and 29 are rejected under 35 U.S.C. 103(a) as being unpatentable over *Sturim* in view of *Waibel*, further in view of *Hermansky* as applied to claims 10 and 25 above, and further in view of *Liu* (6,615,170).

Sturim in view of *Waibel* disclose the method as set forth in claims 10 and 25, however neither disclose applying temporal sequential smoothing to the frame tag using temporal information associated with the clustered anchor model outputs. *Liu* discloses temporal smoothing tagged frames (column 5 line 55- column 5 line 20). *Liu* discloses tagging speech frames based on the output of specific model. Adjacent observations are then used to update the value of a tag for each frame by weighting observations at different times.

Therefore it would have been obvious to one of ordinary skill in the art at the time of the invention to apply temporal sequential smoothing to the frame tags in *Sturim* and *Waibel*, since it enables the incorporation of adjacent frame tags for updating and validating a current frame tag, thus increasing tagging accuracy, as indicated in *Liu* (column 5 lines 64-65).

Conclusion

Any inquiry concerning this communication or earlier communications from the examiner should be directed to Matthew Baker whose telephone number is (571)270-1856. The examiner can normally be reached on 4-5-9, First Friday Off.

If attempts to reach the examiner by telephone are unsuccessful, the examiner's supervisor, Richemond Dorvil can be reached on (571)272-7602. The fax phone number for the organization where this application or proceeding is assigned is 571-273-8300.

Information regarding the status of an application may be obtained from the Patent Application Information Retrieval (PAIR) system. Status information for published applications may be obtained from either Private PAIR or Public PAIR. Status information for unpublished applications is available through Private PAIR only. For more information about the PAIR system, see <http://pair-direct.uspto.gov>. Should you have questions on access to the Private PAIR system, contact the Electronic Business Center (EBC) at 866-217-9197 (toll-free). If you would like assistance from a USPTO Customer Service Representative or access to the automated information system, call 800-786-9199 (IN USA OR CANADA) or 571-272-1000.

4/22/2009

/Talivaldis Ivars Smits/
Primary Examiner, Art Unit 2626

/Matthew Baker/
Examiner, Art Unit 2626